

RubiX and Locality Aware Scheduling

PRESTO SUMMIT 2019, Bangalore

RubiX

Why Caching

- Popularity of Cloud Stores like S3
 - + Near-infinite capacity
 - + Inexpensive
 - + Ease of use
 - Network Latencies
 - Back-offs

RubiX

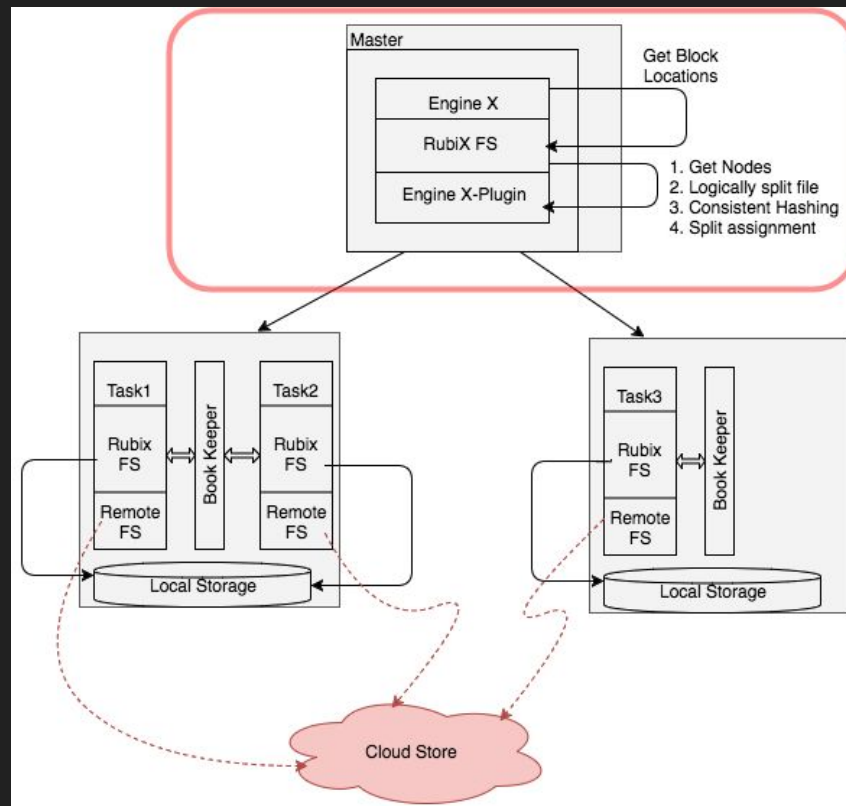
- Extendible to many engines
- Columnar format friendly
- Works well with autoscaling
- Share-able across engines/instances

Architecture

- Split ownership assignment system
- Data Caching System
- Plugins

Architecture

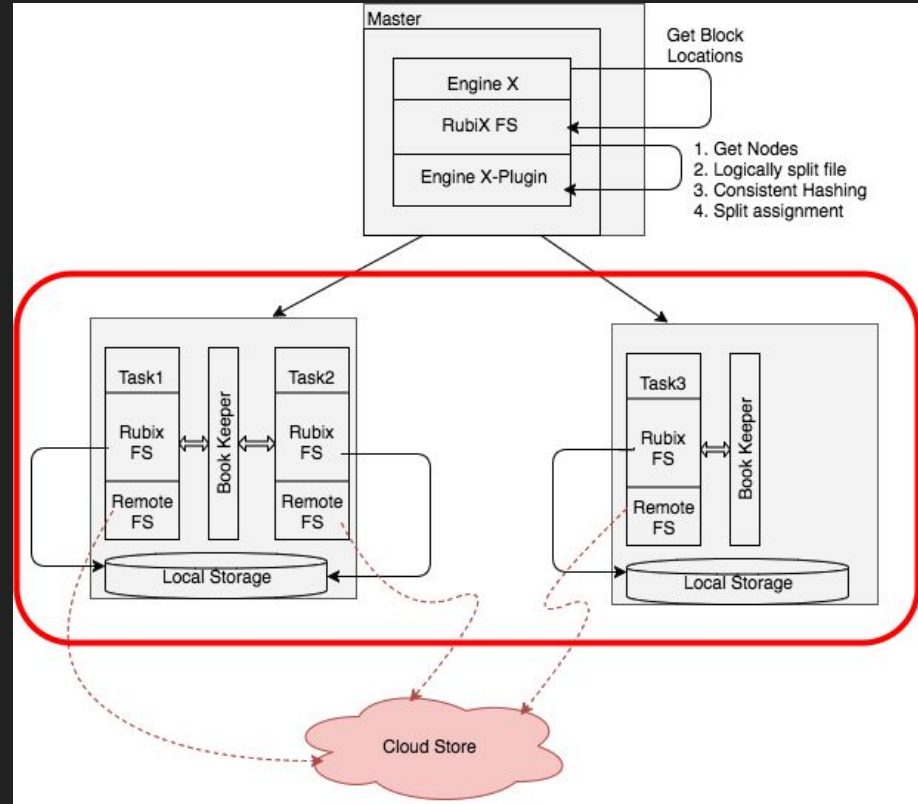
- Split ownership assignment system
 - Used in master node during split computation
 - Calculates which node owns particular split of file
 - Uses Consistent Hashing to work well with Autoscaling



Architecture

- Data Caching System

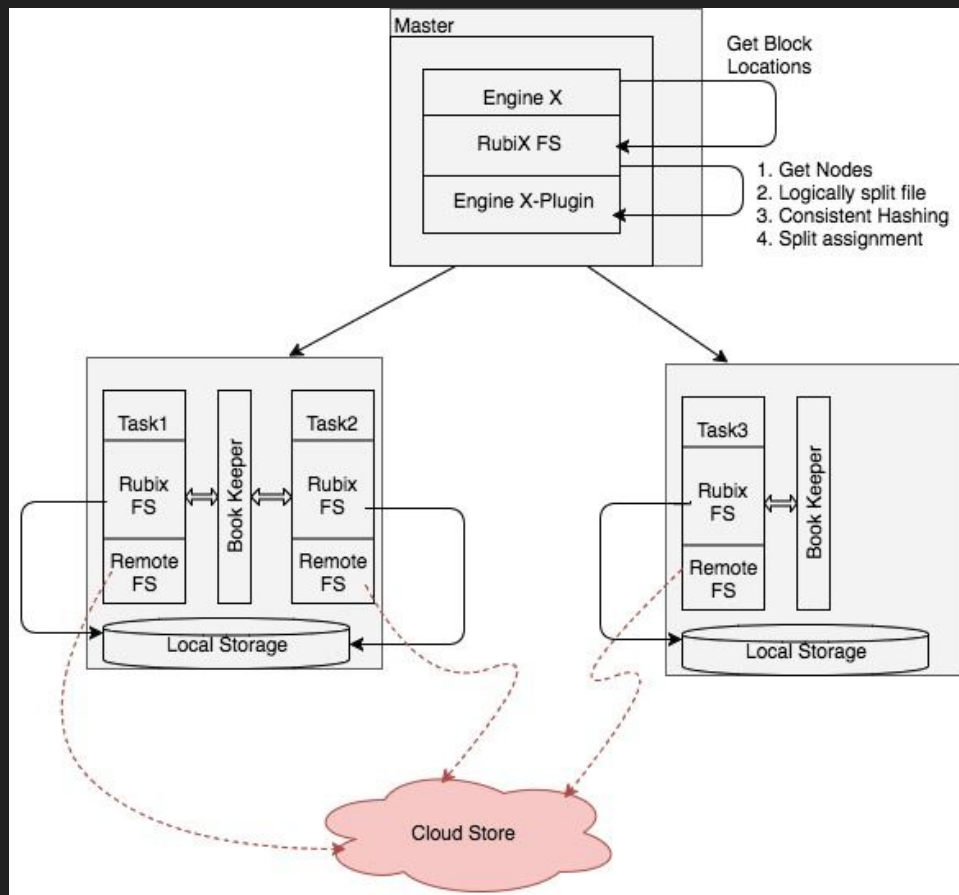
- Used in worker nodes when data is read
- Read from disk or remote as per the metadata
- Metadata stored in units of block (1MB each)
- **BookKeeper** provides metadata for the block
- Metadata too Checkpointed to local disk



Architecture

- Plugin

- Provides two types of information
 - How to get the list of nodes in the system
 - FileSystem for remote reads
- E.g. presto plugin, hadoop1 plugin, hadoop2 plugin



Locality Aware Scheduling in Presto

Presto Scheduler

- Split assignment to nodes
- Degree of parallelism
- Uniform distribution

Presto Scheduler: Advantages

- * Uniform division of workload
- Dynamic: Splits occur in batches

Scope for improvement?

- Uniform division of workload => Uniform query execution time?
- Data locality: What if the data resides within the disk or in a cache?
- Impact on query performance and network traffic?

Locality Aware Scheduling

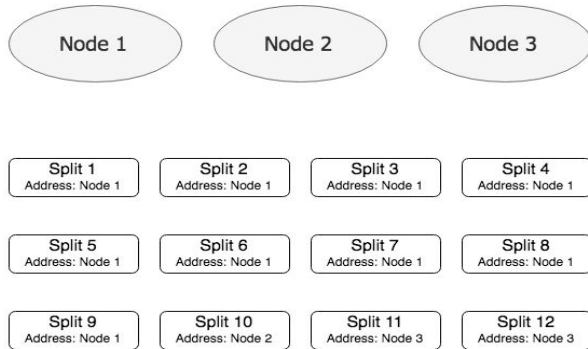
- Locality as an extra dimension for Split assignment
- Ensure uniform workload
- Multi-stage best effort implementation

Locality Aware Scheduling: Implementation

- Three Stages:
 - Stage 1: Assignment based on locality
 - Stage 2: Compute based assignment: Allocate splits to nodes with max slots available
(existing behaviour)
 - Stage 3: Redistribution of splits if non-uniform assignment

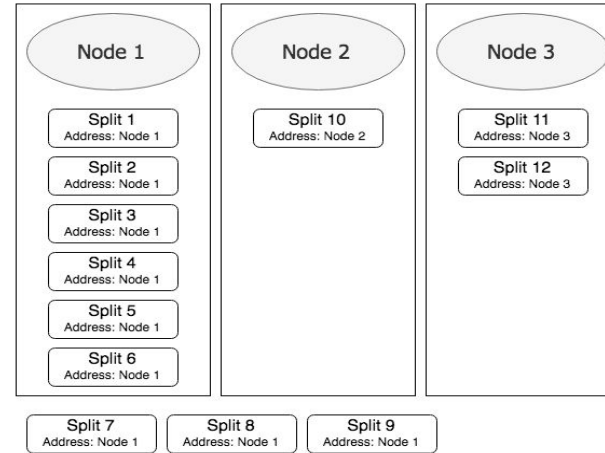
12 Splits to be assigned to 3 Nodes

Max Splits Per Node: 6



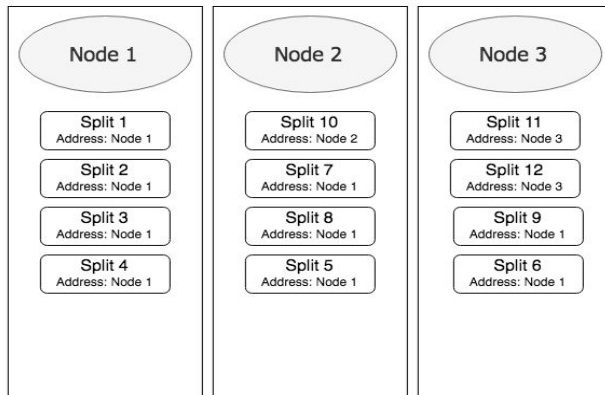
Stage 1 : Locality Based Assignment

Max Splits Per Node: 6



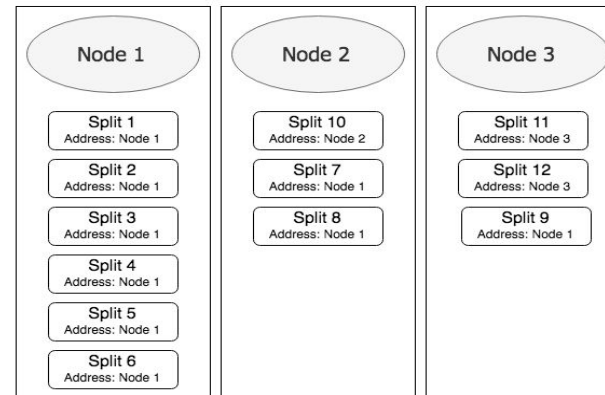
Stage 3 : Redistribution for Uniform Workload

Max Splits Per Node: 6



Stage 2 : Compute Based Assignment

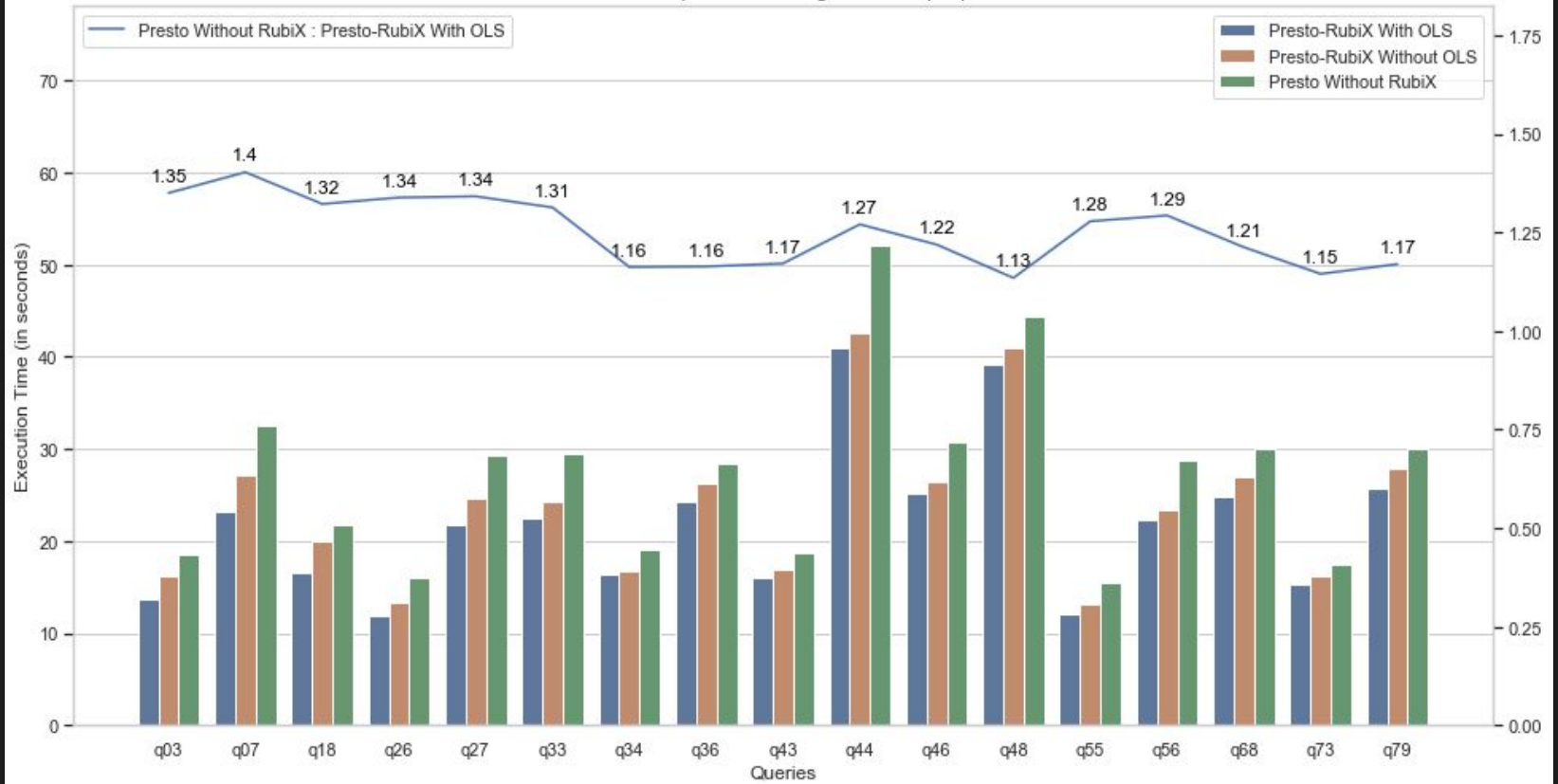
Max Splits Per Node: 6



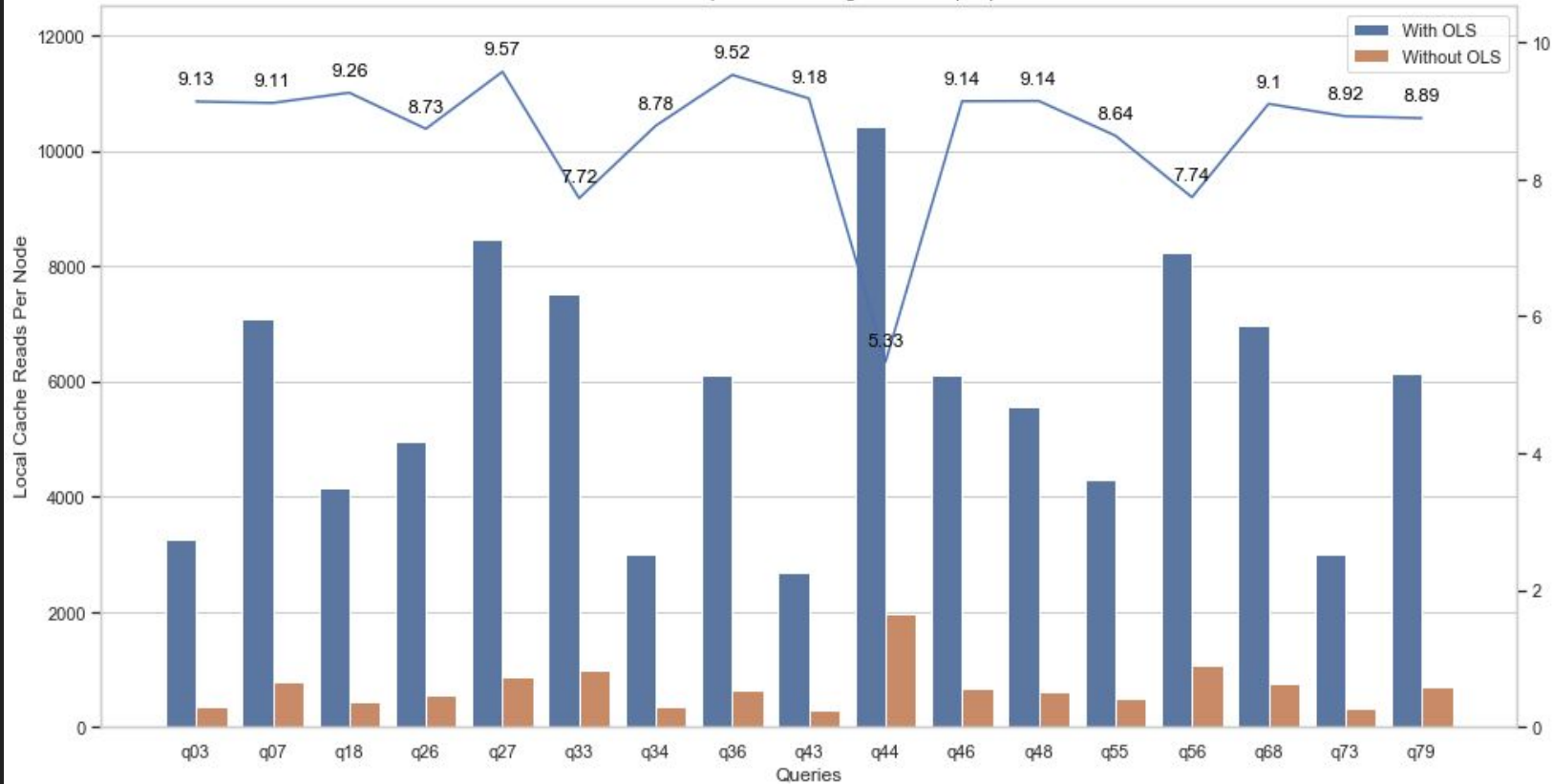
Locality Aware Scheduling: Benefits

- Massive increase in the number of cached/local reads (around 9x)
- Reduced network traffic
- Improved Query perf

Performance comparison i3.4xLarge 10 nodes parquet1000



Cache Reads comparison i3.4xLarge 10 nodes parquet1000



Locality Aware Scheduling

- Has been adopted as the default scheduling model in Presto (release 315)
- Additional scheduler complexity: in the order of milliseconds

Thank you